1    **Family-specific genotype arrays increase the accuracy of pedigree based**

2    **imputation at very low marker densities**

3    Andrew Whalen, Gregor Gorjanc, and John M Hickey

4

5    The Roslin Institute and Royal (Dick) School of Veterinary Studies, The University of

6    Edinburgh, Midlothian, Scotland, UK

7

8                                    **Abstract**

9          In this paper we evaluate the performance of using a family-specific low-density genotype

10    arrays to increase the accuracy of pedigree based imputation. Genotype imputation is a widely

11    used tool that decreases the costs of genotyping a population by genotyping the majority of

12    individuals using a low-density array and using statistical regularities between the low-density and

13    high-density individuals to fill in the missing genotypes. Previous work on population based

14    imputation has found that it is possible to increase the accuracy of imputation by maximizing the

15    number of informative markers on an array. In the context of pedigree based imputation, where

16    the informativeness of a marker depends only on the genotypes of an individual's parents, it may

17    be beneficial to select the markers on each low-density array on a family-by-family basis. In this

18    paper we examined four family-specific low-density marker selection strategies, and evaluated

19    their performance in the context of a real pig breeding dataset. We found that family-specific or

20    sire-specific arrays could increase imputation accuracy by 0.11 at 1 marker per chromosome, by

21    0.027 at 25 markers per chromosome and by 0.007 at 100 markers per chromosome. These results

22    suggest that there may be a room to use family-specific genotyping for very-low-density arrays

23    particularly if a given sire or sire-dam pairing have a large number of offspring.

24

25  **Introduction**

26       In this paper we evaluate the performance of using family-specific low-density genotyping

27  arrays for pedigree based imputation. The use of genomic information in livestock breeding has

28  risen substantially over the past decade, and has led to an increase in the accuracy of selection,

29  particularly on traits with low heritability (Van Eenennaam et al., 2014), decreased the

30  generational interval for some species (notably cattle; Wiggans et al., 2017), and increased the rate

31  of genetic gain (Knol et al., 2016). Many of these gains have been made possible due to the use of

32  low-cost genotypes obtained through genotype imputation. In the context of an animal or plant

33  breeding program, genotype imputation allows most of the individuals in the population to be

34  genotyped with a low-cost, low-density genotype array, while only a small number of individuals

35  (e.g., the sires and top dams) are genotyped with a high-density array. The markers on the low-

36  density array are used to identify shared haplotypes between low-density and high-density

37  individuals. The shared haplotype segments are then used to fill-in missing genotypes (Li and

38  Stephens, 2003).

39       High imputation accuracy is key for maximizing the rate of genetic gain in a population;

40  low imputation accuracy decreases genomic prediction accuracy, which in turn decreases the

41  response to selection. One of the primary factors that influences imputation accuracy is the total

42  number of markers on a low-density genotyping array. If there are too few markers then it may be

43  challenging to correctly identify the shared haplotypes between low-density and high-density

44  individuals. Having more markers increases the specificity of detecting shared haplotypes.

45  However, increasing the number of low-density markers also increases the cost, potentially

46  limiting the number of individuals genotyped. An alternative way to increase accuracy is to keep

47    the total number of markers constant, but choose the markers to be as informative as possible

48    (Aliloo et al., 2018; Boichard et al., 2012; Wu et al., 2016).

49        Past work on population based imputation has found that selecting markers that have high

50    minor allele frequency, are evenly spaced throughout the chromosome (Wu et al., 2016), or covary

51    strongly with other markers can improve imputation accuracy (Aliloo et al., 2018). These three

52    factors allow a population based imputation method to distinguish between high-density reference

53    haplotypes and find the specific reference haplotype that the low-density individual carries. For

54    example, markers with high minor allele frequency are likely to segregate between haplotypes,

55    allowing similar haplotypes to be distinguished. In contrast, markers with a low minor allele

56    frequency may be fixed in most of the haplotypes in the population and so provide limited

57    information.

58        We can also search for informative markers in the context of pedigree based imputation.

59    Unlike for population based imputation methods, where we need to identify which haplotype an

60    individual carries from all of the haplotypes in the population, in pedigree based imputation where

61    an individual is imputed based on the genotypes of their parents, we only need to identify which

62    parental haplotypes the individual inherited at each marker. This reduces the number of haplotypes

63    that need to be considered from hundreds or thousands to just four (for diploid species).

64    Informative markers are those that allow us to distinguish between the parental haplotypes. If the

65    parents have high-density genotypes (potentially by having been imputed themselves) and are

66    phased then the informative markers will be the markers that are heterozygous in the parents. To

67    illustrate this, suppose there is a biallelic marker where both parents are genotyped and phased. If

68    sire is AB and the dam is BB, then the marker is informative for distinguishing sire haplotypes.

69    The resulting child will either be AB or BB. If the child is AB we know that the child inherited the

70    A allele from the sire and since the sire is phased we know which haplotype the child inherited at

71    that marker. Alternatively if the child is BB, we know it inherited the B allele from the sire and

72    the corresponding haplotype. If both parents are heterozygous at a marker (AB and AB), then the

73    marker will be informative for both parents in half of the time, i.e. when the child is either AA or

74    BB. If the child is AB the marker will not be informative since we cannot determine the parent of

75    origin for each allele. We illustrate these conditions in Figure 1.

76         The fact that the marker informativeness for pedigree based imputation is based only on

77    the genotypes of the sire and dam of an individual suggests selecting the markers on the low-

78    density array on a family-by-family basis, by targeting markers that are heterozygous in one, or

79    both parents. In this paper we use simulation to evaluate the performance of four family-specific

80    low-density marker selection strategies and three population based strategies.  In each simulations

81    we used a marker selection strategy to construct a series of low-density arrays. These arrays were

82    then used to mask high-density genotype data taken from a commercial pig population. We used

83    multi-locus iterative peeling (Whalen et al., 2017) to re-impute each individual to high-density.

84    We found that although family-specific genotyping arrays greatly increased the accuracy of

85    imputation at very low marker densities (5-10% gains at < 25 markers per chromosome) but that

86    the gains at low-density arrays with more markers were small (<1%, at >100 markers per

87    chromosome).

88    **Materials and Methods**

89    **Genetic data**

90    In this study, we used genotypes for 1,000 focal individuals and their ancestors from a large

91    commercial pig breeding program. The focal individuals were selected to have been genotyped on

92    a high-density array (~50k markers across 18 chromosomes), and to have had 5 generations of

93 (potentially low-density) genotyped ancestors. In total, we extracted the genotypes for 2,405

94 animals (1,000 focal individuals and 1,405 ancestors). We have then performed several simulations

95 where the genotypes of the focal individuals were masked according to a low-density marker

96 selection strategy (explained below) and imputed using AlphaPeel. AlphaPeel is a pedigree based

97 imputation method based on multi-locus peeling (Whalen et al., 2017;

98 https://alphagenes.roslin.ed.ac.uk/wp/software/alphapeel/). We have run AlphaPeel with default

99 parameters.

100 **Marker selection strategies**

101 We evaluated two sets of marker selection strategies where the markers on the low-density

102 array were either optimized for the whole population, or for a specific family. For all methods, we

103 split the chromosome into $k$ bins, where $k$ is the number of low-density markers, and used a marker

104 selection strategy to select a marker from each bin. For each marker selection strategy, we varied

105 the number of low-density markers per chromosome between 1 and 700 in 16 increments, using

106 either 1,2, 3, 5, 10, 15, 25, 50, 100, 150, 200, 300, 400, 500, 600, or 700 markers.

107 We evaluated three population based marker selection strategies. We selected either the

108 middle marker from each bin (*midpoint*), the marker in the bin that had the highest minor allele

109 frequency (*maf*), or the marker that was simultaneously central and had a high minor allele

110 frequency (*combined*). The combined centrality and minor allele frequency was based on the

111 method of Wu et al. (2016). For each marker we calculated a score:

112 $$score_i = -(1 - d_i)(p_i \log_2(p_i) + (1 - p_i) \log_2(1 - p_i)),$$

113 where $d_i$ is the distance (in number of markers) between the marker and the center of the bin, and

114 $p_i$ is the minor allele frequency for marker $i$. The term $(1-d_i)$ gives higher weight to markers that

115 are close to the center of the bin. The term $(p_i \log_2(p_i) + (1 - p_i) \log_2(1 - p_i))$ is the Shannon

116    information content of the marker based on the minor allele frequency and is highest for markers

117    with minor allele frequency close to 0.5 (Wu et al. 2016). Unlike Wu et al. (2016) we did not

118    perform a global optimization of the location of each markers, but instead selected the marker for

119    each bin independently.

120        Previous work has found that selecting two markers from the first and last bins on the

121    chromosome can improve accuracy (Boichard et al., 2012) due to the higher-than normal

122    recombination rate at the ends of the chromosome. Due to the small number of markers used in

123    our study (in some cases, only 1 marker was used) we only selected a single marker from each bin,

124    even for the first and last bins.

125        We evaluated four family-specific marker selection strategies. We selected the marker

126    closest to the center of the bin that was either heterozygous in both parents (*Het/Het*), heterozygous

127    in one parent and homozygous in the other (*Het/Hom*), heterozygous in at least one parent

128    (*Het/Any*), or heterozygous in the sire (Het/Sire). In the Het/Hom condition we used $\frac{k}{2}$ bins and

129    separately selected markers in each bin that were informative for the sire or the dam (if the number

130    of markers was odd, the sire received $\frac{k+1}{2}$ bins, and the dam received $\frac{k-1}{2}$ bins). If a bin did not

131    contain an acceptable marker for the family-specific strategy, we used the *combined* population

132    strategy to select the marker for that bin. This occurred primarily in the Het/Het condition when

133    there were no markers that were heterozygous in both parents, or when the number of low-density

134    markers was large. In all family-specific strategies, we restricted the pool of potential markers to

135    markers that were genotyped in the real dataset (i.e., not missing) in the sire, dam, and offspring.

136    **Imputation accuracy measurement**

137        Imputation accuracy was measured as the correlation between an individual's imputed

138    genotype and their true genotype, corrected for their parent average genotype:

139       $accuracy = cor(G_{imputed} - G_{parent\_average},\ G_{true} - G_{parent\_average})$.

140      This measure of imputation accuracy is designed specifically for pedigree based imputation. It is

141      0 if no genotype information is available on a focal individual (leading the individual to be imputed

142      as the parent average genotype), and is 1 if the individual is imputed perfectly. The goal of this

143      metric is to assess the accuracy of imputing within-family (Mendelian sampling) genotype

144      variation. In simulations we have found a close relationship between this measure of imputation

145      accuracy and the accuracy of the breeding value estimates. In addition, this measure does not rely

146      on using the population minor allele frequency (as opposed to correcting for minor allele

147      frequency, as in Calus et al., 2014), which may not be representative of the allele frequencies in

148      specific families. In cases where the genotypes of the parents were missing in the real dataset, we

149      used the imputed values from AlphaPeel to calculate the parent average genotype. This was

150      primarily done to fill in spontaneous missing genotypes, and to impute dams that were genotyped

151      at a lower density.

152           Imputation accuracies were calculated separately for each chromosome and then averaged

153      across all 18 chromosomes.

154      **Results**

155           In Figure 2, we present the performance of using either a population strategy or a family-

156      specific strategy, for both the (a) absolute imputation accuracy, or (b) imputation accuracy relative

157      to the *combined* population strategy. We found that the *combined* strategy was the highest

158      performing population strategy, followed by the *maf* strategy, and then by the *midpoint* strategy.

159      The difference between the *combined* strategy and the *maf* strategy was less than 0.001 at above

160      25 markers per chromosome. Of the family-specific strategies, we found that the *Het/Hom* strategy

161      was the highest performing strategy, followed by the *Het/Any* strategy, and finally the *Het/Het*

162 strategy. The *Het/Sire* strategy performed better than the Het/Het strategy with fewer than 5

163 markers, but worse with 5 or more markers. For all marker densities, the family-specific strategies

164 outperformed the *combined* strategy.

165 The *combined* strategy gave high imputation accuracies across a range of marker densities.

166 Imputation accuracy was 0.312 at 1 marker per chromosome (18 markers total), 0.796 at 10

167 markers per chromosome (180 markers total), 0.903 at 25 markers per chromosome (450 markers

168 total), 0.945 at 50 markers per chromosome (900 markers total), and 0.985 at 500 markers per

169 chromosome (9,000 markers genome wide).

170 Using a family-specific strategy further increased imputation accuracy. When the Het/Any

171 strategy was used, we obtained an 0.111 gain in imputation accuracy compared to the *combined*

172 strategy at 1 marker per chromosome. This dropped to 0.058 at 10 markers per chromosome, 0.027

173 at 25 markers per chromosome, 0.014 at 50 markers per chromosome, and 0.010 at 500 markers

174 per chromosome. The gains for the other family-specific strategies were similar.

175 In Figure 3(a), we plot the imputation accuracy with the Het/Any strategy by chromosome,

176 and in Figure 3(b) by chromosome length. We found that imputation accuracy decreased as the

177 chromosome length increased, but that this difference was small even for large chromosomes. To

178 quantify these differences in imputation accuracy, we used a linear model to measure the effect of

179 the number of markers and chromosome length (in cM) on accuracy. Chromosome lengths were

180 taken from Tortereau et al. (2012). The linear model fitted chromosome length as a linear covariate

181 nested within the number of markers as a categorical variable to account for the non-linear effect

182 that number of markers has on accuracy. We found a significant effect of chromosome length on

183 accuracy (regression coefficients decreased from -0.0012 loss of accuracy per cM at 2 marker per

184    chromosome to 0.0001 loss of accuracy per cM at 100 marker per chromosome, p<0.001) and the

185    interaction between the number of markers and chromosome length (p<0.001).

186    **Discussion**

187          In this paper we evaluate the performance of using family-specific low-density marker

188    selection strategies to increase the accuracy of pedigree based imputation. We found that using

189    parental genotype information to select markers on a low-density genotype array could increase

190    imputation accuracy, with the largest gains occurring at very low marker densities (between a 0.11

191    and 0.05 increase in accuracy for between 1 and 25 markers per chromosome). The gains were

192    more limited at higher marker densities (under a 0.01 increase in accuracy at more than 100

193    markers per chromosome). In addition, we quantified the influence that chromosome length had

194    on imputation accuracy, and found that increasing chromosome length had a near-linear impact on

195    imputation accuracy when the number of informative markers per chromosome was held constant.

196    In the remainder of the discussion we will highlight the performance of each family-specific

197    marker selection strategy, compare our results to past work on optimizing the design of low-density

198    arrays for population based imputation, and discuss the commercial viability of using family-

199    specific genotype arrays.

200

201    **Performance of family-specific  marker selection strategies**

202          In this paper we found that selecting the markers on a low-density genotype array based on

203    parental information increased accuracy in all cases compared to using the same set of markers for

204    every individual in the population. We evaluated four marker selection strategies, and found that

205    selecting markers that were heterozygous in one parent, and homozygous in the other (Het/Hom,

206    Figure 1a) yielded the highest imputation accuracies. Selecting markers that were heterozygous in

207    both parents (Het/Het, Figure 1b) resulted in much lower imputation accuracies than the Het/Hom

208    strategy, particularly at very low marker densities. This effect is caused by the lack of informative

209    markers when low-density individuals have heterozygous genotypes in the Het/Het condition

210    (Figure 1b).

211         In addition to the strategies presented in Figure 1, we also investigated two hybrid

212    strategies. The first was to select markers that were heterozygous in either (Het/Any). The

213    second was to select markers that were heterozygous in the sire (Het/Sire). We found that the

214    Het/Any strategy performed in between the Het/Hom and Het/Het strategies, reflecting the fact

215    that markers chosen were split between being heterozygous in one parent and homozygous in the

216    other, and being heterozygous in both parents. We found that the Het/Sire condition performed

217    well at a few markers per chromosome, but that the gain in imputation accuracy declined more

218    rapidly compared to the alternative strategies. This is likely due to the Het/Sire strategy placing

219    most of its weight on finding markers that are informative for the sire, resulting in few markers

220    that were informative for the dam. Even so, the Het/Sire strategy outperformed all of the

221    population strategies tested, making it a potentially useful strategy when a single sire produces a

222    large number of offspring.

223         One of the advantages of studying family-specific marker selection strategies is that

224    because they focus all of their genotyping efforts on informative loci, they also provide anupper

225    bound on the performance any population-specific strategy. We found that the difference between

226    all of the family-specific strategies and the worst performing population strategy was less than

227    0.01 at 100 markers, suggesting that for pedigree based imputation there are limited gains for

228    optimizing the design for low-density arrays if more than 100 markers per chromosome are used

229    (1,800 markers in total for the 18 pig autosomal chromosomes in our study population).

230

**Comparison to population based imputation**

The results in this paper align closely with the previous work on optimizing low-density genotyping arrays for population based imputation. Similar to both Aliloo et al. (2018) and Wu et al. (2016), we find that the gains in imputation accuracy for an optimized array were highest at low marker densities and diminished at higher densities. We were also able to replicate the primary finding of Wu et al. (2016), that simultaneously optimizing the low-density markers for both high minor allele frequency and even spacing improved imputation accuracy particularly at low-densities.

Consistent with past work on population and pedigree based imputation (Antolín et al., 2017) we found that the accuracy of pedigree based imputation was higher than that of population based imputation at similar marker densities. This is expected because population based imputation has to compare an individual's low-density genotype to all of the population haplotypes, while pedigree based imputation has to match it to the four parental haplotypes. When the number of low-density markers is small it is hard to distinguish among population haplotypes, but much easier to distinguish among parental haplotypes. When the number of markers increases, distinguishing population haplotypes becomes easier. Therefore, in the context of optimizing the low-density arrays, family-specific strategies will be relevant only at low marker densities. For example, Aliloo et al. (2018) obtained a gain in imputation accuracy of 0.10 at ~125 markers per chromosome using an optimized set of markers and a population based imputation algorithm (absolute imputation accuracy rose from 0.69 to 0.79). In contrast, we observed an accuracy of 0.97 at 100 markers per chromosome with pedigree based imputation, obtained a gain in imputation accuracy of 0.10 at 3 markers per chromosome (going from 0.55 to 0.65 accuracy) in the Het/Any condition.

253

**Commercial feasibility of family based imputation**

254

255        The primary question of using family-specific genotype arrays revolves around the cost

256     and the complexity of deploying such arrays in the context of a genetic improvement program.

257     There are two primary issues: First, in order for a family-specific array to be beneficial, the density

258     of the array needs to be low. Second, the use of a family-specific array may require the construction

259     of a large number of arrays, which may be prohibitively expensive. We discuss both of these issues

260     in more detail below.

261        On the question of marker densities, we find that in order for a family-specific genotype

262     array to be beneficial, the underlying marker density has to be much smaller than what is traditional

263     used in an animal improvement program (<25 markers per chromosome), and will result in lower

264     absolute values of imputation accuracy than a traditional low or medium density array. This limits

265     the use case for family-specific arrays into the situation where having imperfect genetic

266     information is acceptable, i.e., to cases where the accuracy of selection can be low, or when

267     selection decisions are not directly made on the genotyped individual. Such situations might

268     include genotyping individuals in a non-nucleus environment to establish flow of phenotypic

269     information to individuals in the nucleus, or performing genetic improvement in breeding

270     programs where very low-density arrays are used to genotype a very large number of offspring.

271     This might have potential in aquaculture (Lillehammer et al., 2013; Tsai et al., 2017) and  crop

272     breeding (Gonen et al., 2018; Jacobson et al., 2015).

273        On the question of the number of arrays, because the family-specific genotype arrays

274     depend on the genotypes of both the sire and the dam, the number of different arrays individuals

275     in the population need to be genotyped at may be large. This will be particularly the case when a

276    single dam has a limited number of offspring (most notably in cattle and small ruminants, but also

277    in pigs). In these cases it may be possible to reduce the number of arrays needed by using a sire-

278    specific genotype array. Alternatively, there may be situations where a single sire-dam pair may

279    produce a large number of offspring as is the case in aquaculture and crop breeding, or where a

280    more flexible genotyping method could be deployed (Thomson et al., 2012).

281    **Conclusion**

282    Overall this paper evaluates the utility of family information to select markers on a low-

283    density array. Although we find minimal gains at the densities currently used in modern breeding

284    programs (over 100 markers per chromosome), we find high increases in accuracy at very low

285    marker densities (between 1-25 markers per chromosome), and may be particularly useful when

286    expanding genotyping efforts to individuals that are not traditionally genotyped.

287

288    **Acknowledgements**

294

295    **References**

296    Aliloo, H., Mrode, R., Okeyo, A.M., Ojango, J., Dessie, T., Rege, J.E.O., Goddard, M.E.,

297    and Gibson, J.P. (2018). Optimal design of low-density marker panels for genotype imputation.

298    Proc. Fo World Congr. Genet. Appl. Livest. Prod.

299        Browning, S.R., and Browning, B.L. (2007). Rapid and accurate haplotype phasing and

300    missing-data inference for whole-genome association studies by use of localized haplotype

301    clustering. Am. J. Hum. Genet. *81*, 1084–1097.

302        Das, S., Forer, L., Schönherr, S., Sidore, C., Locke, A.E., Kwong, A., Vrieze, S.I., Chew,

303    E.Y., Levy, S., and McGue, M. (2016). Next-generation genotype imputation service and

304    methods. Nat. Genet. *48*, 1284–1287.

305        Hickey, J.M., Kinghorn, B.P., Tier, B., Wilson, J.F., Dunstan, N., and Werf, J.H. van der

306    (2011). A combined long-range phasing and long haplotype imputation method to impute phase

307    for SNP genotypes. Genet. Sel. Evol. *43*, 12.

308        Knol, E.F., Nielsen, B., and Knap, P.W. (2016). Genomic selection in commercial pig

309    breeding. Anim. Front. *6*, 15.

310        Meuwissen, T., and Goddard, M. (2010). The Use of Family Relationships and Linkage

311    Disequilibrium to Impute Phase and Missing Genotypes in Up to Whole-Genome Sequence

312    Density Genotypic Data. Genetics *185*, 1441–1449.

313        Sargolzaei, M., Chesnais, J.P., and Schenkel, F.S. (2011). FImpute - An efficient

314    imputation algorithm for dairy cattle populations. J. Dairy Sci. *94 (E-Suppl. 1)*, 421.

315        Van Eenennaam, A.L., Weigel, K.A., Young, A.E., Cleveland, M.A., and Dekkers,

316    J.C.M. (2014). Applied Animal Genomics: Results from the Field. Annu. Rev. Anim. Biosci. *2*,

317    105–139.

318        Whalen, A., Ros-Freixedes, R., Wilson, D.L., Gorjanc, G., and Hickey, J.M. (2017).

319    Hybrid peeling for fast and accurate calling, phasing, and imputation with sequence data of any

320    coverage in pedigrees. bioRxiv.

321        Wiggans, G.R., Cole, J.B., Hubbard, S.M., and Sonstegard, T.S. (2017). Genomic
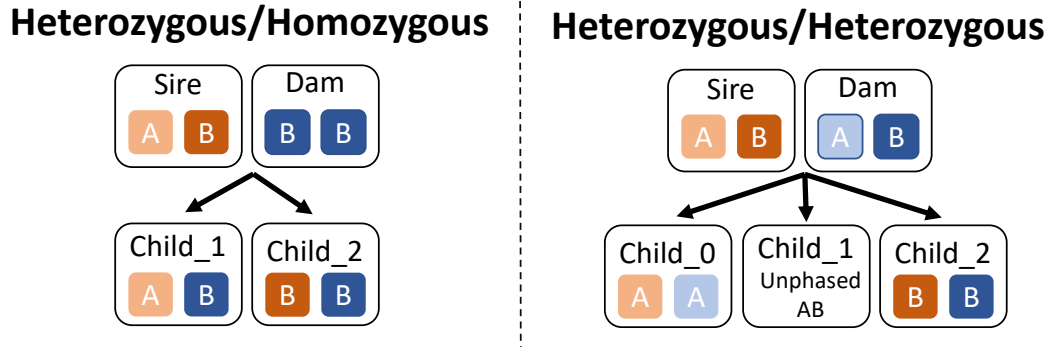
322    Selection in Dairy Cattle: The USDA Experience. Annu. Rev. Anim. Biosci. *5*, 309–327.

323          Wu, X.-L., Xu, J., Feng, G., Wiggans, G.R., Taylor, J.F., He, J., Qian, C., Qiu, J.,

324    Simpson, B., Walker, J., et al. (2016). Optimal Design of Low-Density SNP Arrays for Genomic

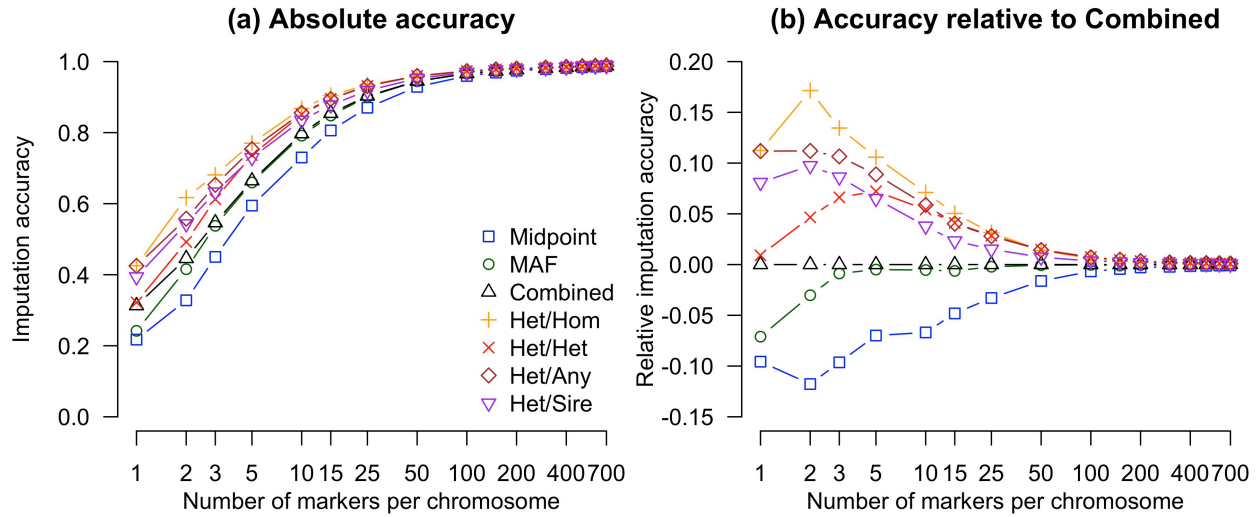325    Prediction: Algorithm and Applications. PLOS ONE *11*, e0161719.

326

327

328            Figure 1. A graphical representation of informative markers for pedigree based

329     imputation.

330

Figure 2. Imputation accuracy as a function of the number of markers per chromosome and the marker selection strategy. P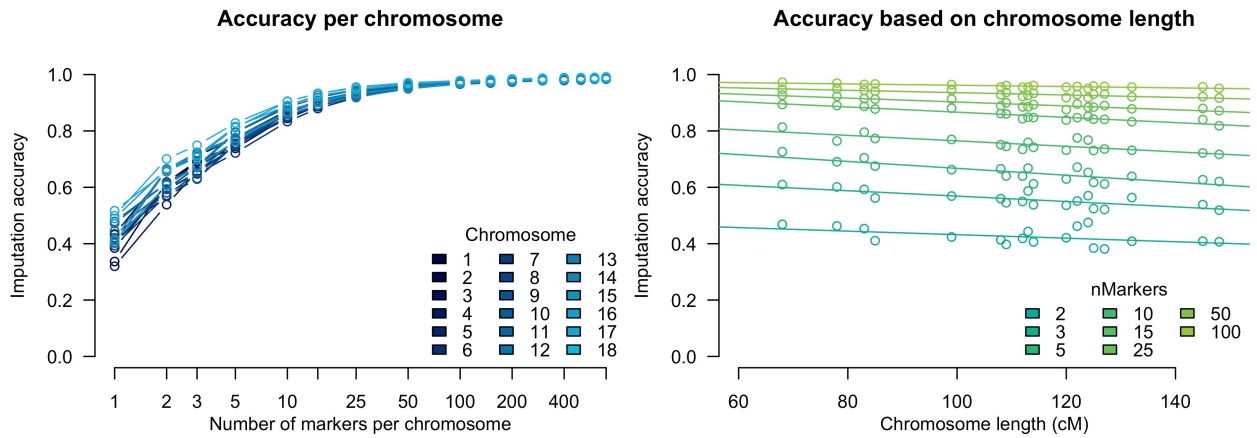anel (a) provides the absolute imputation accuracy (measured as correlation between the true and imputed genotypes of an individuals corrected for parent average genotype), while panel (b) provides comparison relative to the *combined* strategy.

337

338          Figure 3. Imputation accuracy by (a) chromosome and (b) chromosome length. In both

339   panels the Het/Any strategy was used to select the markers on the low-density arrays.

340